



STATISTICS FOR THE DATA SCIENCE

Part - 3

- POISSON DISTRIBUTION
- NORMAL / GAUSSIAN DISTRIBUTION
- UNIFORM DISTRIBUTION
- Z - SCORE
- CENTRAL LIMIT THEOREM
- ESTIMATOR
- HYPOTHESIS AND TESTING MECHANISM
- P - VALE
- Z - TEST

#Value_freeContent



@Krishan kumar

$$* \text{Pmf} = P(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$$

for $k = 0, 1, 2, 3, \dots, n$ where

$${}^n C_k = \frac{n!}{k! (n-k)!}$$

\Rightarrow Mean :-

$$\text{Mean} = np$$

\Rightarrow Variance :-

$$\text{Var} = npq$$

\Rightarrow Std

$$\text{Std} = \sqrt{npq}$$

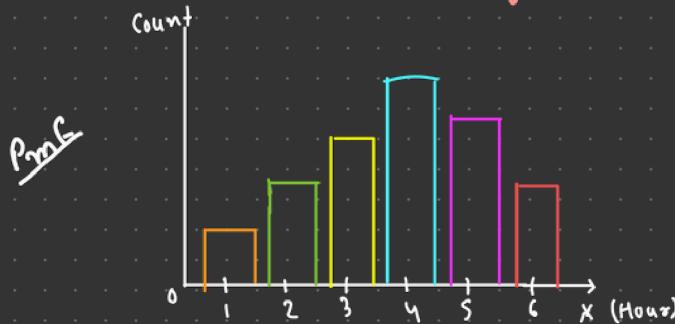
► Poisson Distribution

* Discrete random variable. (Pmf)

* Describe the number of event occurring in fixed time interval.

Ex) ① No. of people visiting hospital every hour:

② No. of people visiting bank every hour.



$\lambda=3$, Expected no. of event occur at every time interval.

Means How many people come at every hour.

Q What is the probability of a person to come at 5th hour.

Ans

$$\text{Pmf} = P(X=5) = \frac{e^{\lambda} - \lambda^X}{L^X}, \text{ if } \lambda = 3 \quad (\text{Ans})$$
$$= \frac{e^3 - 3^5}{L^5} \Rightarrow 0.101 \rightarrow 10.1\%$$

there is 10% possibilities that at 5th hour 3 person come.

Q What is the probability to visit at 5th hour OR 4th hour.

Ans
Pmf

$$P(X=5) + P(X=4)$$

* Mean of poission distribution

$$\text{Mean} = E(X) = \mu$$

$$\boxed{\mu = \lambda \times t}$$

λ = expected no. of event
occur at every time
interval

t = time interval

$$\boxed{\text{Variance of poission} = \lambda \times t}$$

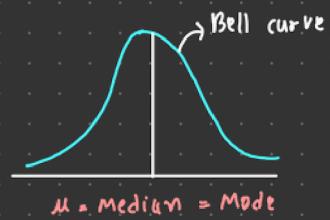
► Normal / Gaussian Distribution \rightarrow (PdF)

Notation :- $N(\mu, \sigma^2)$

Parameter :- $\mu \in \mathbb{R}$ ^{Real number} (Mean)

$\sigma^2 \in \mathbb{R} > 0$ = variance

$x \in \mathbb{R}$ = Data points



$$\text{PDF} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{\sigma^2}}$$

⇒ Mean of normal distribution

Mean = μ = Average

⇒ Variance

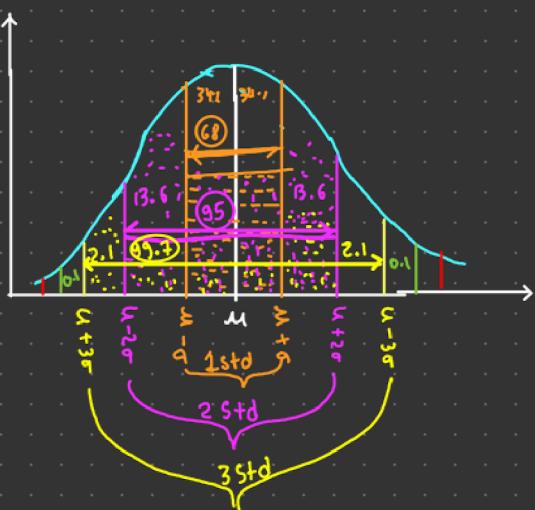
$$\text{Var} = \sigma^2$$

This rule follow

68-95-99.7% Rule

⇒ Std

$$\sigma = \sqrt{\text{Var}}$$



$$P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

- Ex,
- ① Weight of the students in the class.
 - ② Height of " " " "
 - ③ IRIS Dataset [Sepal width]

► Uniform Distribution →

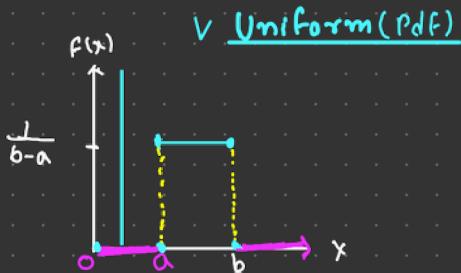
(A) Continuous uniform distribution (Pdf)

(B) Discrete Uniform distribution (Pmf)

(A) Continuous Uniform Distribution,

In statistic, the continuous Uniform dis. OR Rectangular dis. is a family of symmetric probability dis. The dis. describe an experiment where there is an arbitrary outcome that lie between certain bounds.

the bounds are defined by the parameter a and b, which are the minimum and maximum va

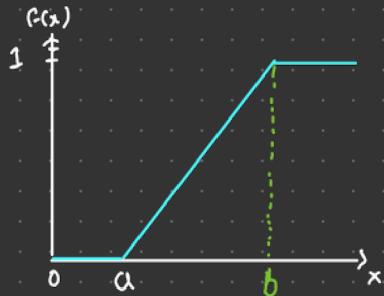


Notation → $V(a,b)$

Parameter,

$$-\infty < a < b < \infty$$

► Cdf



$$p.d.F \rightarrow \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & , \text{ otherwise} \end{cases}$$

$$c.d.F \rightarrow \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a, b] \\ 1 & \text{for } x > b \end{cases}$$

$$\text{Mean} = \frac{1}{2}(a+b)$$

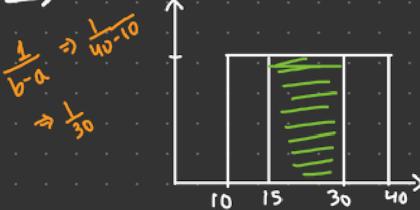
$$\text{Variance} \rightarrow \frac{1}{12}(b-a)^2$$

$$\text{Median} \rightarrow \frac{1}{2}(a+b)$$

Eg The no. of candies sold daily at a shop is uniformly dis. with a maxi at 40 and mini of 10.

① Probability of daily sales falls between 15 and 30.

Sol:



$$\lambda_1 = 15$$

$$\lambda_2 = 30$$

$$P(15 \leq X \leq 30) = (\lambda_2 - \lambda_1) \times \frac{1}{b-a}$$

$$\Rightarrow 15 \times \frac{1}{30} \Rightarrow \frac{1}{2} \Rightarrow 0.5$$

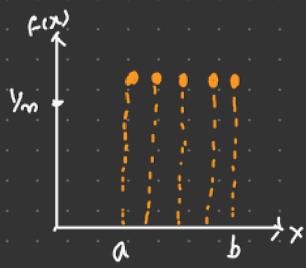
② $P(X \geq 20)$

Sol. $(40-20) \times \frac{1}{30}$

$\Rightarrow 0.66 \rightarrow 66\%$

(B) Discrete uniform dis. (Pmf)

In statistic the discrete uniform dis. is symmetric probability dis. wherein a finite number of value are equally likely to be observed every one of n value has equally probability $\frac{1}{n}$. Another way of saying that "discrete uniform Dis" would be a known finite number of outcome equally likely to happen.



e.g Rolling a dice

$$[1, 2, 3, 4, 5, 6]$$

$$\begin{aligned} \hookrightarrow P(1) &\rightarrow \frac{1}{6} & a &\rightarrow 1 \\ &\vdots && \\ P(6) &\rightarrow \frac{1}{6} & b &\rightarrow 6 \end{aligned}$$

Notation $\rightarrow M(a, b)$

$$b-a+1 = n$$

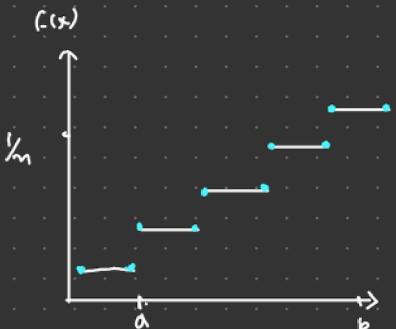
Parameter,

a, b with $b \geq a$

$$\text{Pmf} = \frac{1}{n}$$

$$\text{Mean} = \frac{a+b}{2}$$

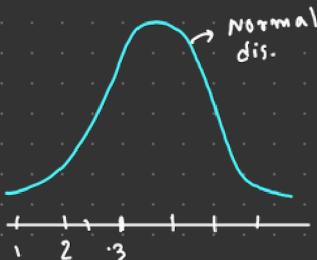
Median



Standard Normal Distribution and Z-score

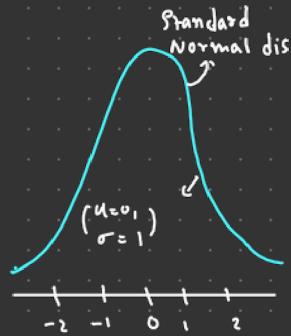
$$X \in [1, 2, 3, 4, 5], \mu = 3$$

$$\sigma = 1.414 \approx 1$$



Transition Technique

When
 $\mu = 0$
 $\sigma = 1$



\Rightarrow Z-Score

Z-score is used for the help for Normal dis. to transform in standard normal distribution.

$$\boxed{\text{Z-Score} = \frac{x_i - \mu}{\sigma}}$$

$$x_i = 1, \frac{1-3}{1} \Rightarrow -2$$

$$x_i = 4 \Rightarrow \frac{4-3}{1} \rightarrow 1$$

$$x_i = 2, \frac{2-3}{1} \Rightarrow -1$$

$$x_i = 3 \Rightarrow \frac{3-3}{1} \rightarrow 0$$

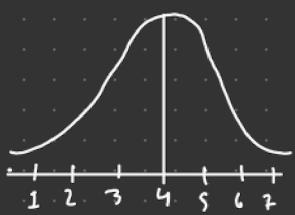
$$x_i = 3, \frac{3-3}{1} \Rightarrow 0$$

Q Away from the mean for a specific mean how much Std fall for the 4.

Ans by Z-score, $x_i = 4 \Rightarrow \frac{4-3}{1} \rightarrow 1$

Note :

One Std to the right \rightarrow When value is +ve
One Std to the left \rightarrow When value is -ve



$$\left\{ \begin{array}{l} \mu = 4 \\ \sigma = 1 \end{array} \right\} \quad (\text{Normal dis.fun.})$$

Q How many standard deviation 4.5 is away from the mean?

Ans $x_i = 4.5$

$$\text{Z-score} \rightarrow \frac{4.5 - 4}{1} \rightarrow 0.5$$

0.5 Std from the right.

Q What percentage of data is falling above 4.5



$$\mu \rightarrow 4$$

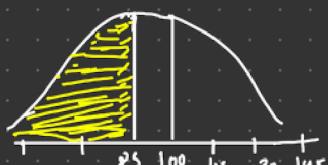
$$\sigma \rightarrow 1$$

$$\text{Z-score} \rightarrow \frac{4.5 - 4}{1} \Rightarrow 2.5$$

Area under the curve (≤ 2.5) $\Rightarrow 1 - 0.6681 \rightarrow 6.6\%$

Problem In India average IQ is 100, with a std of 15. What is the percentage of the population whose you expect to have an IQ lower than 85.

$$\mu \rightarrow 100, \sigma = 15$$

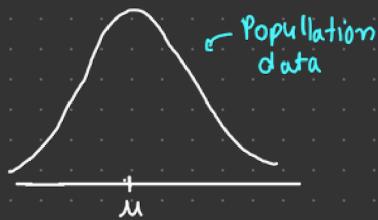


$$\text{Z-score} \rightarrow \frac{85 - 100}{15} \rightarrow -\frac{15}{15} \rightarrow -1$$

Area under the curve (≤ 85) $= 0.15866 \rightarrow 15.8\%$

Central limit theorem

type, ① $X \sim N(\mu, \sigma)$



$n = 20$
sample distribution

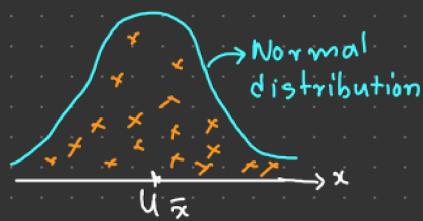
$$S_1 \rightarrow [x_1, x_2, \dots, x_{10}] = \bar{x}_1$$

$$S_2 \rightarrow [x_2, x_3, \dots, x_{20}] = \bar{x}_2$$

⋮

$$S_n \rightarrow \bar{x}_n \Rightarrow \text{Sample mean.}$$

$$\bar{x} = \{\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_n\}$$

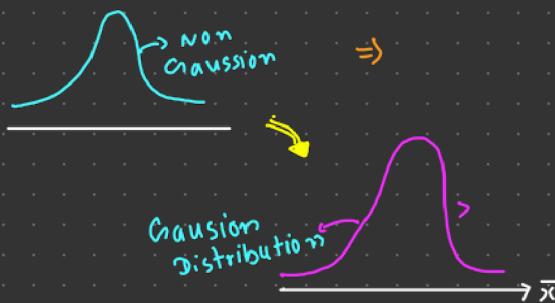


If we have population data who is normally distributed. On this we applied Sampling distri. and after that calculating sample mean.

↳ when we plotted all those mean \rightarrow We got normal distribution

② $X \not\sim N(\mu, \sigma)$

$n \geq 30$

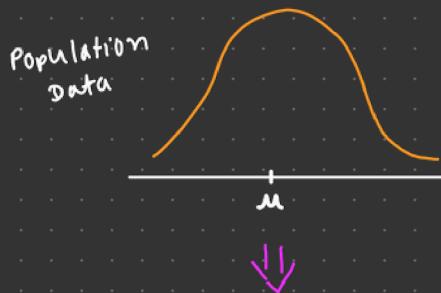


$S_1 \rightarrow \{x_1, x_2, \dots, x_{30}\} \rightarrow \bar{x}_1$
 $S_2 \rightarrow \{x_2, x_3, \dots, x_{30}\} \rightarrow \bar{x}_2$
⋮
 $S_n \rightarrow \bar{x}_n$

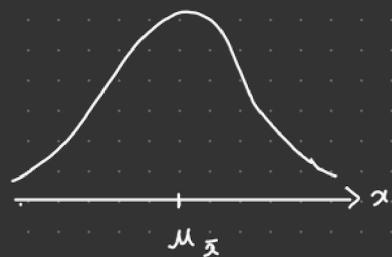
⇒ CLT Properties

the central limit theorem says that the sampling dis. of the mean will be always be normally distributed, as long as the sample size is large enough. Regardless of whether the population has a normal, poisson, binomial or any other dist. the sampling distribution of the mean will be normal.

Important for the interview



Sampling distribution
of mean ((LT))



$$x \sim N(\mu, \sigma)$$

σ → population std

μ → Population mean

n → Sample size.

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

n can be any value

Q What is the Std in sample distribution?

S

$$\text{Std} \Rightarrow \frac{\sigma}{\sqrt{n}}$$

Infrential Statistics

Estimate:

It is an observe numerical value used to estimate an unknown population parameter.

① Point estimate,

Single numerical value used to estimate the unknown population parameter.

* Sample mean is a point estimate of a population mean.



There is a huge gape, to counter those gape we use Interval estimate.

② Interval estimate,

Range of value used to estimate the unknown population parameter.

* Interval estimate of population parameter are called confidence interval.



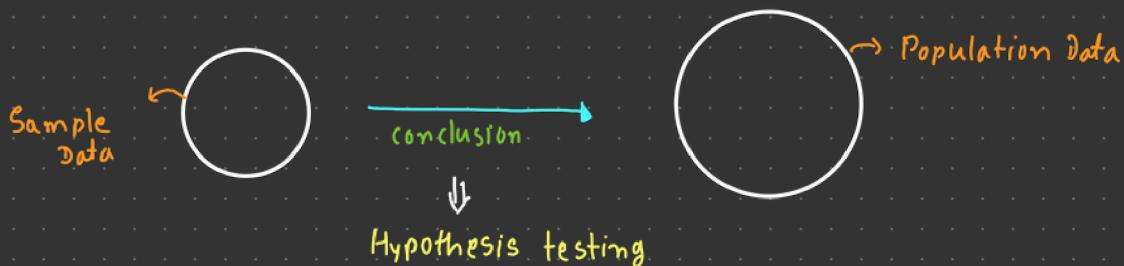
55 - 65





A Hypothesis and Hypothesis testing mechanism:-

↳ In inferential Stats. used for Conclusion or Inferences.



* Hypothesis testing Mechanism

↔ Person ~~not~~ done crime

(1) Null Hypothesis (H_0) → The person is not Guilty
↳ The assumption that you are beginning with.

(2) Alternate Hypothesis (H_1), The person is Guilty

↳ Opposite of Null Hypothesis

(3) Experiments → prove collect (by help of \rightarrow Statistical Analysis)
↳ { DNA, Tests, Finge prints } \rightarrow (p-value)

(4) Accept the null hypothesis OR Reject the null hypothesis

for example,

College at district a stats says its average passed percentage of statistic are 85%. A New college opened in the district and it was found that a sample of 100 Students have a pass % of 90 with a Standard deviation of 40%. Does this school have a different Passed%.

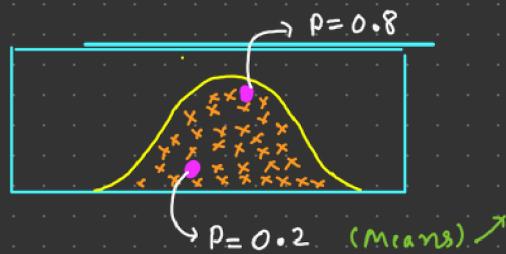
Ans → Null Hypothesis $\rightarrow (H_0) \rightarrow \mu = 85\%$

Alternate hypothesis $\rightarrow (H_1) \rightarrow \mu \neq 85\%$

P-Value

The P-value is a number, calculated from a statistical test, that describes how likely you are to have found a particular set of observations if the null hypothesis were true. P-values are used in hypothesis testing to help decide whether to reject the null hypothesis.

Ex) Using keyboards spare top.



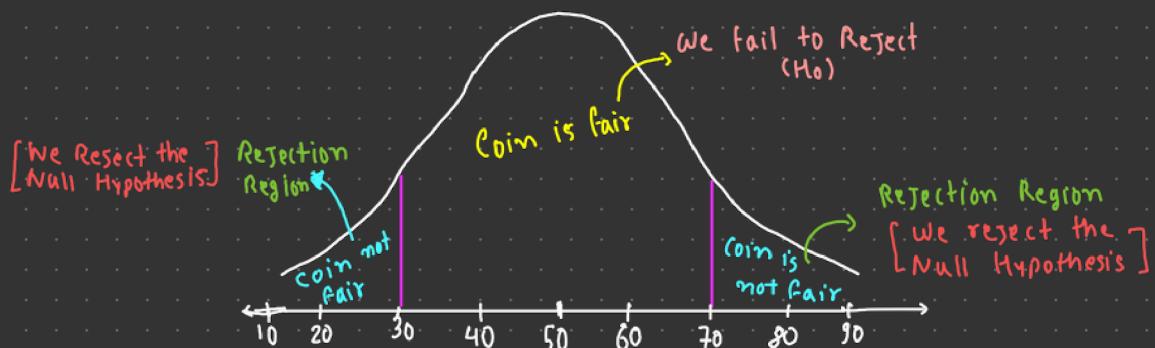
out of 100 touches in this key the probability of touching in this region ≥ 0.8

↳ Hypothesis testing

Ex) Whether coin is fair or Not. f 100 Tosses

Ans)

- ① Null Hypothesis (H_0): coin is fair
- ② Alternate hypothesis (H_1): Coin is not fair
- ③ Experiment;



* Significance value (α)

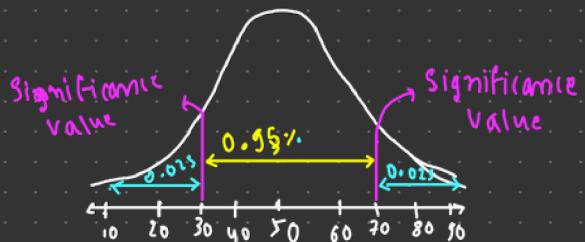
Let,

$$\alpha = 0.05$$

confidence interval (C.I.)

$$C.I. \Rightarrow 1 - 0.05$$

$$\Rightarrow \underline{\underline{0.95}}$$



* Conclusion

Let, $P = 0.01\%$

$\Rightarrow P < \text{Significance}$

\therefore We reject the null hypothesis.

else:

We fail to reject null hypothesis.



Hypothesis testing and statistical Analysis:-

- (A) Z - Test } Average
- (B) T - Test }
- (C) CHI SQUARE } categorical data
- (D) ANNOVA } variance

(A) Z-test :-

Condition : Z-score only apply on where,

- (i) Population std
- (ii) $n \geq 30$

Problem

The average height of all the residents in a city is 168cm. with a $\sigma = 3.9$. A Doctor believe that mean to be different. He measure the height of 36 individuals and found the average height to be 169.5 cm.

(a) State Null and alternate hypothesis.

(b) At 95% confidence level, is there enough evidence to reject the null hypothesis.

Method 1 \rightarrow Z-test

$$\mu = 168\text{cm}, \sigma = 3.9, n = 36, \bar{x} = 169.5\text{cm}$$

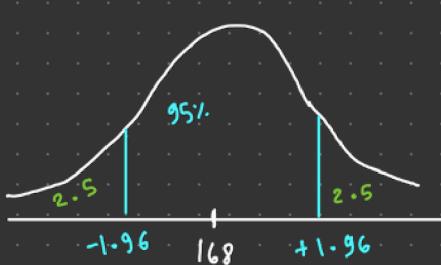
(a) Null hypothesis (H_0) = $\mu = 168\text{cm}$

Alternate hypothesis (H_1) = $\mu \neq 168$ {2 Tail test}

It can be greater OR
lesser than also.

$$(b) C.I. = 0.95, \alpha = 1 - 0.95 \\ = \underline{\underline{0.05}}$$

Decision boundary



* Area under the curve

$$1 - 0.95 \Rightarrow 0.05 \\ 0.25 \leftarrow \overbrace{\quad}^{0.25}$$

$$\Rightarrow 1 - 0.25 \Rightarrow \boxed{0.9750}$$

Area under the curve

$$Z\text{-Score} \rightarrow 0.9750 \rightarrow +1.96$$

Z-Score for population data

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

Z-Score for sample data

$$Z\text{-Score} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

X Statistical Analysis

$$Z\text{-Score} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \Rightarrow \frac{169.5 - 168}{3.9 / \sqrt{36}} \Rightarrow 2.31$$

If Z-test value is less than -1.96 OR greater than +1.96 we

Reject the Null hypothesis

else

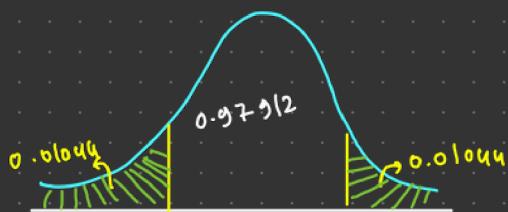
we Accept the Null hypothesis

$2.31 > +1.96$ { So, we Reject the null hypothesis }

Method 2 → p-value

* Statistical Analysis

$$Z\text{-Score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}} \Rightarrow 2.31$$



In Z-table $2.31 \Rightarrow 0.98956$

$$\begin{aligned} & 1 - \text{Area under the curve} \\ & \Rightarrow 1 - 0.98956 \\ & \Rightarrow \boxed{0.01044} \end{aligned}$$

* How to calculate p-value

$$\Rightarrow 0.01044 + 0.01044$$

$$\Rightarrow \underline{\underline{0.02088}}$$

If P value < Significance

$$0.02088 < 0.05$$

{We Rejecting the null hypothesis}

else-

we accept the null hypothesis.

↳ Here we accepting the null hypothesis.

Problem: 2

A factory manufacture bulbs with a average warrenty of 5 yrs with standard deviation of 0.50%. A worker believe that the bulb will manufacture in less then 5 yrs. He test a sample of 40 bulbs and find the average time to be 4.8 Years.

(a) Stats the null hypothesis and alternate

(b) At a 2% Significance level , is there enough evidence to support the idea that the warrenty should be revised.

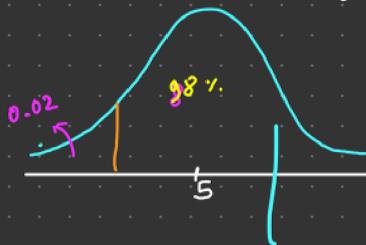
$$\rightarrow \mu = 5, \sigma = 0.50, \bar{x} = 4.8, n = 40$$

\hookrightarrow (A) Null hypothesis (H_0) $\Rightarrow \boxed{\mu = 5}$

\hookrightarrow Alternate hypothesis (H_1) $\Rightarrow \underline{\mu < 5}$ { 2 Tail test }
 (only checking one side)

* Method (P-Value)

- Decision Boundary

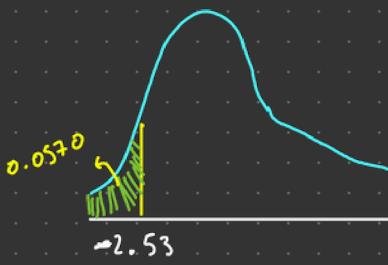


C.I. $\rightarrow 0.98$

$$\alpha = 1 - 0.98 \\ \Rightarrow \boxed{0.02}$$

\rightarrow p-value \rightarrow

$$z\text{-test} \rightarrow \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \rightarrow \frac{4.8 - 5}{0.05/\sqrt{40}} \rightarrow \frac{-0.2}{0.079} \rightarrow -2.53$$



\hookrightarrow Area under the curve of
 -2.53 , z value is $= \boxed{0.0570}$

$$p\text{-value} = 0.0570$$

if $p\text{value} < \text{Significance}$

$$0.0570 < 0.02 \rightarrow \text{False}$$

* Conclusion

The warranty needs to be revised. So, we accept the Null hypothesis.

